

Optimistic Active Learning using Mutual Information

Yuhong Guo and Russell Greiner

ALL AND ALL AN

University of Alberta, Edmonton, Canada

Active Learning: A process of sequentially deciding which unlabeled instance to label, with the goal of producing the best classifier with limited number of labelled instances.

Idea: an optimistic active learner that exploits the discriminative partition information in the unlabeled instances, makes an optimistic assessment of each candidate instance, and temporarily switches to a different policy if the optimistic assessment is wrong.

Optimistic Query Selection

1. Most uncertain query selection (MU):

$$\begin{split} & \underset{i \in U}{\operatorname{argmax}} H(\left.Y_{i} \left| \left.\mathbf{x}_{i}, L\right.\right) \\ & H(\left.Y_{i} \left| \left.\mathbf{x}_{i}, L\right.\right) = -\sum_{y_{i}} P(\left.y_{i} \left| \left.\mathbf{x}_{i}, L\right.\right) \log P(\left.y_{i} \left| \left.\mathbf{x}_{i}, L\right.\right) \right. \end{split}$$

Shortcoming: ignores the unlabeled data !

2. Select the query that maximizes its conditional mutual information about the unlabeled data:

(a) Take the expectation wrt
$$Y_i$$
 (MCMI[avg]):

$$\operatorname{argmin}_i \sum_y P(y | \mathbf{x}_i, \boldsymbol{\theta}_L) H(Y_u | \mathbf{x}_u, \boldsymbol{\theta}_{L+(\mathbf{x}_i, y)})$$

Shortcoming: aggravates the ambiguity caused by the limited labelled data.

(b) Take an optimistic strategy: use only the best query label (MCMI[min]):

 $\underset{i \in U}{\operatorname{argmin}} f(i)$ where $f(i) = \min_{y} \sum_{u} H(Y_u \,|\, \mathbf{x}_u, \boldsymbol{\theta}_{L+(\mathbf{x}_i, y)})$

 $\begin{array}{l} L \quad \Im \quad \text{set of } labeled \text{ instances} \\ U &= \text{ index set for } unlabeled \text{ instances} \\ X_U \quad \Im \quad \text{set of all unlabeled instances} \end{array}$

Online Adjustment

Optimistic query selection tries to identify the well separated partition. E.g.:



Potential problem: given only a few labelled data points, there might be *many* settings leading to well-separated classes

example:



Online Adjustment :

- Can easily detect this "guessed wrong" situation, in the immediate next step,
- Simply compare the actual label for the query with its optimistically predicted label
- Whenever Mm+M guesses wrong,
- it switches to a different query selection criterion (MU) for the next 1 iteration

Mm+M Algorithm



Experimental Evaluation Comparing Mm+M with other Active Learners						• Empirical results
AUSTRALIAN" BREAST" CLEVE" CORRAL CRX DIABETES" FLARE" GERMAN" GLASS2" HEART" HEPATITIS" MOFN	$\begin{array}{c} 587 & 0 & (+) \\ 0 \ / \ 0 & (-) \\ 567 & 0 & (+) \\ 577 & 14 & (0) \\ 807 & 0 & (+) \\ 537 & 0 & (+) \\ 537 & 0 & (+) \\ 387 & 0 & (+) \\ 387 & 0 & (+) \\ 387 & 0 & (+) \\ 377 & 0 & (+) \\ 777 & 0 & (+) \\ 777 & 0 & (+) \\ 777 & 0 & (+) \\ 777 & 0 & (+) \\ 777 & 0 & (+) \\ 777 & 0 & (+) \\ 777 & 0 & (+) \\ 7777 & 0 & (+) \\ 7777 & 0 & (+) \\ 77777 & 0 & (+) \\ 77777 & 0 & (+) \\ 777777 & 0 & (+) \\ 7777777 & 0 & (+) \\ 77777777777777777777777777777777777$	$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	$\begin{array}{c} 11 / 14 (-) \\ 88 / 0 (+) \\ 26 / 0 (+) \\ 13 / 0 (+) \\ 2 / 8 (-) \\ 11 / 12 (+) \\ 27 / 2 (+) \\ 0 / 0 (+) \\ 8 / 0 (+) \\ 0 / 0 (+) \\ 0 / 0 (-) \\ 0 /$	$\begin{array}{c} 167 & 0 & (+) \\ 0 & /53 & (-) \\ 777 & 0 & (+) \\ 0 & /23 & (-) \\ 137 & 0 & (+) \\ 737 & 8 & (+) \\ 597 & 1 & (+) \\ 917 & 0 & (+) \\ 917 & 0 & (+) \\ 297 & 0 & (+) \\ 507 & 0 \\ 507 & 0 \end{array}$	show Mm+M works better than • MU • MCMI[min] • MCMI[avg] • MU (MU-SVM) • Random
PIMA* VOTE	85/ 2 (+) 0/2/(-)	2 / 0 (+) 97 / 0 (+)) 84/ 2 (+)) 0/ 0 (-)	0/0(0) 2/0(+)	36/13 (+) 0/16 (-)	Future work:
IRIS VEHICLE [*] LYMPHOGRAPHY	68/ 1 (+) 56/ 0 (+) 0/ 4 (0)	17/ 0 (+) 84/ 0 (+) 0 / 1 (-)) 13 / 23 (0)) 38 / 0 (+)) 0 / 30 (-)	38/ 1 (+) 0 / 19 (-) 0 / 2 (-)	$ \begin{array}{rrrr} -I & -0 \\ -I & -0 \\ -I & -0 \\ \end{array} $	•Understand when Mm+M is appropriate
TOTAL W/L/T Signed Rank Test	12 / 3/ 2 12/ 3/ 2	10 / 1/6 12/2/3	13/2/2 13/1/3	7 / 2 / 8 10 / 3 / 4	10 / 4 / 0 9 / 4 / 1	•Design further variants
Comparing Mm+M Over 100 sample s • "statistically bette • "statistically wors • "tied" 13 times Signed Rank Test	vs MU, for PIM izes, Mm+M wa r" 85 times e" 2 times "shows" Mm+M	A dataset: C as M is better	Comparing Mm+M vs MCMI[avg], over 17 datasets: Mm+M was "statistically better" for >5 more sample-sizes: 13 times "statistically worse" for >5 more sample-sizes: 2 times Signed Rank Test "shows" Mm+M is better: 13 times			

worse: 1 time